CLAIMS

What is claimed is:

1.    A method for managing data, the method comprising the steps of:

      maintaining a plurality of persistent data items on persistent storage accessible to a plurality of nodes, the persistent data items including a particular data item stored at a particular location on said persistent storage;

      assigning exclusive ownership of each of the persistent data items to one of the nodes, wherein a particular node of said plurality of nodes is assigned exclusive ownership of said particular data item;

      when any node wants an operation performed that involves said particular data item, the node that desires the operation to be performed ships the operation to the particular node for the particular node to perform the operation on the particular data item as said particular data item resides at said particular location;

      while the first node continues to operate, reassigning ownership of the particular data item from the particular node to another node without moving the particular data item from said particular location on said persistent storage;

      after the reassignment, when any node wants an operation performed that involves said particular data item, the node that desires the operation to be performed ships the operation to said other node for the other node to perform the operation on the particular data item as said particular data item resides at said particular location.

-28-

1     2.       The method of Claim 1 wherein the step of reassigning ownership of the particular

2     data item from the particular node to another node includes updating an ownership map that

3     is shared among the plurality of nodes.

1     3.       The method of Claim 1 wherein the plurality of nodes are nodes of a multi-node

2     database system.

1     4.       The method of Claim 3 wherein the multi-node database system includes nodes that

2     do not have access to said persistent storage.

1     5.       The method of Claim 3, wherein:

2          said persistent storage is a first persistent storage of a plurality of persistent storages

3                used by said multi-node database system; and

4          the method further comprises reassigning ownership of a second data item from a first

5                node that has access to said first persistent storage to a second node that has

6                access to a second persistent storage but does not have access to said first

7                persistent storage; and

8          wherein the step of reassigning ownership of the second data item includes moving

9                the second data item from said first persistent storage to said second persistent

10               storage.

1     6.       The method of Claim 3 wherein the step of reassigning ownership of the particular

2     data item from the particular node to another node is performed in response to the addition of

3     said other node to said multi-node database system.

1     7.       The method of Claim 3 wherein:

50277-2277 (OID 2003-123-01)

2  the step of reassigning ownership of the particular data item from the particular node

3     to another node is performed in anticipation of the removal of said particular

4     node from said multi-node database system; and

5   the method further comprises the step of, in anticipation of the removal of said

6     particular node from said multi-node database system, physically moving

7     from said persistent storage to another persistent storage a second data item

8     that is reassigned from said particular node to a node of said multi-node

9     database system that does not have access to said persistent storage.

1 8.  The method of Claim 3 wherein the step of reassigning ownership of the particular

2 data item from the particular node to another node is performed as part of a gradual transfer

3 of ownership from said particular node to one or more other nodes.

1 9.  The method of Claim 8 wherein the gradual transfer is initiated in response to

2 detecting that said particular node is overworked relative to one or more other nodes in said

3 multi-node database system.

1 10.  The method of Claim 9 wherein the gradual transfer is terminated in response to

2 detecting that said particular node is now longer overworked relative to one or more other

3 nodes of said multi-node database system.

1 11.  The method of Claim 3 wherein the step of reassigning ownership of the particular

2 data item from the particular node to another node is performed as part of a gradual transfer

3 of ownership to said other node by one or more other nodes, wherein said gradual transfer is

4 initiated in response to detecting that said other node is underworked relative to one or more

5 other nodes in said multi-node database system.

50277-2277 (OID 2003-123-01)

1    12. The method of Claim 3 further comprising the steps of:

2        after a first node has been removed from the multi-node system, continuing to have a

3           set of data items owned by the first node; and

4        reassigning ownership of data items from the first node to one or more other nodes in

5           response to detecting requests for operations that involve said data items.

1    13. The method of Claim 3 further comprising the steps of:

2        after a first node has been removed from the multi-node system, continuing to have a

3           set of data items owned by the first node; and

4        reassigning ownership of a data item from the first node to a second node in response

5           to detecting that the workload of said second node has fallen below a

6           predetermined threshold.

1    14.    The method of Claim 1 wherein:

2    at the time said particular data item is to be reassigned to said other node, the

3        particular node stores a dirty version of said particular data item in volatile

4        memory; and

5    the step of reassigning ownership of the particular data item from the particular node

6        to another node includes writing said dirty version of said particular data item

7        to said persistent storage.

1    15.    The method of Claim 1 wherein:

2    at the time said particular data item is to be reassigned to said other node, the

3        particular node stores a dirty version of said particular data item in volatile

4        memory; and

50277-2277 (OID 2003-123-01)

5    the step of reassigning ownership of the particular data item from the particular node

6    to another node includes forcing to persistent storage one or more redo records

7    associated with said dirty version, and purging said dirty version from said

8    volatile memory without writing said dirty version of said particular data item

9    to said persistent storage;

10    said other node reconstructs said dirty version by applying said one or more redo

11    records to the version of the particular data item that resides on said persistent

12    storage.

1  16.    The method of Claim 1 wherein:

2    at the time said particular data item is to be reassigned to said other node, the

3    particular node stores a dirty version of said particular data item in volatile

4    memory; and

5    the method further includes the step of transferring the dirty version of said particular

6    data item from volatile memory associated with said particular node to

7    volatile memory associated with said other node.

1  17.    The method of Claim 16 wherein the step of transferring the dirty version is

2  performed proactively by the particular node without the other node requesting the dirty

3  version.

1  18.    The method of Claim 16 wherein the step of transferring the dirty version is

2  performed by the particular node in response to a request for the dirty version from said other

3  node.

1  19.    The method of Claim 1 wherein:

-32-

2  the step of reassigning ownership of the particular data item from the particular node

3    to another node is performed without waiting for a transaction that is

4    modifying the data item to commit;

5  the transaction makes a first set of modifications while the particular data item is

6    owned by the particular node; and

7  the transaction makes a second set of modifications while the particular data item is

8    owned by said other node.

1 20. The method of Claim 19 further comprising rolling back changes made by said

2 transaction by rolling back the second set of modifications based on undo records in an undo

3 log associated with said other node, and rolling back the first set of modifications based on

4 undo records in an undo log associated with said particular node.

1 21. The method of Claim 1 wherein method includes the steps of:

2  the other node receiving a request to update said data item;

3  determining whether the particular node held exclusive-mode or shared-mode access

4    to the data item;

5  if the particular node did not hold exclusive-mode or shared-mode access to the data

6    item, then the other node updating the particular data item without waiting for

7    the particular node to flush any dirty version of the data item, or redo for the

8    dirty version, to persistent storage.

1 22. The method of Claim 1 further comprising the steps of:

2  in response to transferring ownership of said particular data item to said other node,

3    aborting an in-progress operation that involves said particular data item;

-33-

4      after ownership of the particular data item has been transferred to said particular

5          node, re-executing the in-progress operation.

1   23.    The method of Claim 1 wherein:

2      an operation that involves said particular data item is in-progress at the time the

3          transfer of ownership of said particular data item is to be performed;

4      the method further includes the step of determining whether to wait for said in-

5          progress operation to complete based on a set of one or more factors; and

6      if it is determined to not wait for said in-progress operation to complete, aborting said

7          in-progress operation.

1   24.    The method of Claim 23 wherein said set of one of more factors includes how much

2  work has already been performed by said in-progress operation.

1   25.    A method of managing data, the method comprising the steps of:

2      maintaining a plurality of persistent data items on persistent storage accessible to a

3          plurality of nodes;

4      assigning ownership of each of the persistent data items to one of the nodes by

5          assigning each data item to one of a plurality of buckets; and

6          assigning each bucket to one of the plurality of nodes;

7          wherein the node to which a bucket is assigned is established to be owner of

8             all data items assigned to the bucket;

9      when a first node wants an operation performed that involves a data item owned by a

10          second node, the first node ships the operation to the second node for the

11          second node to perform the operation.

1    26.    The method of Claim 25 wherein the step of assigning each data item to one of a

2    plurality of buckets is performed by applying a hash function to a name associated with each

3    data item.

1    27.    The method of Claim 25 wherein the step of assigning each bucket to one of the

2    plurality of nodes is performed by applying a hash function to an identifier associated with

3    each bucket.

1    28.    The method of Claim 25 wherein the step of assigning each data item to one of a

2    plurality of buckets is performed using range-based partitioning.

1    29.    The method of Claim 25 wherein the step of assigning each bucket to one of the

2    plurality of nodes is performed using range-based partitioning.

1    30.    The method of Claim 25 wherein the step of assigning each data item to one of a

2    plurality of buckets is performed by enumerating individual data-item-to-bucket

3    relationships.

1    31.    The method of Claim 25 wherein the step of assigning each bucket to one of the

2    plurality of nodes is performed by enumerating individual bucket-to-node relationships.

1    32.    The method of Claim 25 wherein the number of buckets is greater than the number of

2    nodes, and the bucket-to-node relationship is a many-to-one relationship.

1    33.    The method of Claim 25 further comprising the step of reassigning from a first node

2    to a second node ownership of all data items that are mapped to a bucket by modifying a

3    bucket-to-node mapping without modifying a data-item-to-bucket mapping.

1

50277-2277 (OID 2003-123-01)

1   34.    A computer-readable medium carrying one or more sequences of instructions which,

2  when executed by one or more processors, causes the one or more processors to perform the

3  method recited in Claim 1.

1   35.    A computer-readable medium carrying one or more sequences of instructions which,

2  when executed by one or more processors, causes the one or more processors to perform the

3  method recited in Claim 2.

1   36.    A computer-readable medium carrying one or more sequences of instructions which,

2  when executed by one or more processors, causes the one or more processors to perform the

3  method recited in Claim 3.

1   37.    A computer-readable medium carrying one or more sequences of instructions which,

2  when executed by one or more processors, causes the one or more processors to perform the

3  method recited in Claim 4.

1   38.    A computer-readable medium carrying one or more sequences of instructions which,

2  when executed by one or more processors, causes the one or more processors to perform the

3  method recited in Claim 5.

1   39.    A computer-readable medium carrying one or more sequences of instructions which,

2  when executed by one or more processors, causes the one or more processors to perform the

3  method recited in Claim 6.

1   40.    A computer-readable medium carrying one or more sequences of instructions which,

2  when executed by one or more processors, causes the one or more processors to perform the

3  method recited in Claim 7.

-36-

1    41.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 8.

1    42.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 9.

1    43.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 10.

1    44.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 11.

1    45.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 12.

1    46.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 13.

1    47.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 14.

1    48.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 15.

1    49.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 16.

1    50.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 17.

1    51.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 18.

1    52.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 19.

1    53.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 20.

1    54.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 21.

1    55.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 22.

1    56.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 23.

1    57.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 24.

1    58.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 25.

1    59.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 26.

1    60.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 27.

1    61.    A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 28.

1    62.     A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 29.

1    63.     A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 30.

1    64.     A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 31.

1    65.     A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 32.

1    66.     A computer-readable medium carrying one or more sequences of instructions which,

2    when executed by one or more processors, causes the one or more processors to perform the

3    method recited in Claim 33.

1    67.     A method for use in a multi-node shared-nothing database system, the method

2    comprising the steps of:

3        a first node of said multi-node shared-nothing database system initially functioning as

4           exclusive owner of a first data item and a second data item, wherein said first

5           data item and said second data item are persistently stored data items within a

6           database managed by the multi-node shared-nothing database system;

50277-2277 (OID 2003-123-01)

| | |
|---|---|
| 7 | without changing the location of a first data item on persistent storage or shutting |
| 8 | down said first node, reassigning ownership of the first data item from the first |
| 9 | node to a second node of said multi-node shared-nothing database system; and |
| 10 | after reassigning ownership, the first node continuing to operate as the owner of the |
| 11 | second data item, and to handle all requests for operations on said second data |
| 12 | item. |

50277-2277 (OID 2003-123-01)